# HIV Database Workshop
## www.hiv.lanl.gov
## seq-info@lanl.gov

**Presenters**: Bette Korber, Brian Foley & Will FIscher

**Database PIs**: Bette Korber, Thomas Leitner, Karina Yusim

**Additional database staff**: Werner Abfalterer, Peter Hraber, Elisabeth Sharon Fung, Robert Funkhouser, Kumkum Ganguly, Jenni Macke, James Szinger, and Hyejin Yoon

**Project Officer**: Stuart Shapiro, NIAID, NIH

*Theoretical Biology and Biophysics, T-6*
*Los Alamos National Laboratory*

---

# Workshop Topics

### HIV Sequence Database and Immunology Database

Bette Korber, Brian Foley and Will Fischer

*Session 1*

Wednesday,
March 12
11:15 – 12:45

*General introduction*
*Sequence search interface – alignments and basic trees*
*Geography search interface*
*Histogram*
*Database Alignments*

*Tools:*
- *Genecutter - processing nucleotide sequences*
- *Neighbor Joining Treemaker*
- *HIV/SIV sequence locator tool*
- *New HIV Gene Map JBROWSE tool*
- *Highlighter*
- *Protein Feature Accent*
- *Quality Control (if time permits)*

# Workshop Topics

**HIV Sequence Database and Immunology Database**

Bette Korber, Will Fischer and Brian Foley

**Session 2**

Thursday,
March 13
11:15 – 12:45

*Immunology database introduction*
*Epitope maps and epitope summary tables*
*T-cell epitope search*
*T-cell epitope variants*
*Antibody search*
*List of most broadly neutralizing antibodies*
*HIV/SIV sequence locator tool*
*QuickAlign – Align an epitope to the database alignments*
*CATNAP*
*ELF – epitope location finder*
*Peptgen – Design peptides for reagent development*

*Mosaic Vaccine Maker, Epicover, and Posicover*
  *- generate candidate vaccines*
  *- estimate epitope coverage*
  *- determine regional epitope coverage*

HIV Database workshop

**Los Alamos**
NATIONAL LABORATORY

---

# Workshop Goals

■ Understanding the database content, how information was obtained, and what is available

■ Database searching

■ Examples of tools for analyses

■ Quality control tools

**Los Alamos**
NATIONAL LABORATORY

## HIV DATABASES

## Entry page at http://www.hiv.lanl.gov/

The **HIV databases** contain data on HIV genetic sequences, immunological epitopes, drug resistance-associated mutations, and vaccine trials. The website also gives access to a large number of tools that can be used to analyze these data. This project is funded by the Division of AIDS of the National Institute of Allergy and Infectious Diseases (NIAID), a part of the National Institutes of Health (NIH). Click on any of the links below to access a database. Editorial Board

**SEQUENCE DATABASE ▶**    **VACCINE DATABASE ▶**

**IMMUNOLOGY DATABASE ▶**    **OTHER VIRUSES ▶**

### News:    Archived News ▶

**New Features for Epitope Location Finder (ELF)**
ELF displays known and predicted epitopes found within a protein sequence query. ELF results now include both Class I (CTL) and Class II (helper) epitopes. In addition to predicting epitopes based on anchor residues, ELF now includes predictions from the Class I and Class II Binding Predictions tools at the Immune Epitope Database (IEDB). *13 March 2012*

**New Features for HIV BLAST**
HIV BLAST has new features. It now allows the user to find best matches among only subtyped sequences, or sequences of a specific subtype. It allows the resulting sequences to be downloaded fully aligned. *01 March 2012*

**New Option for N-GlycoSite**
The N-GlycoSite tool predicts N-linked glycosylation sites in amino acid sequences. A new option allows the user to exclude sites with a second-position proline, which is disfavored for N-linked glycosylation. *29 February 2012*

**HIV Antibody Search Results More Specific**
The antibody search interface in the HIV Immunology database is now more specific. Searches from the Author, Keyword, and Note fields now display only those notes and references that relate directly to the search. The user may still opt to display all, if desired *09 February 2012*

**New Options for Quickalign**
The Quickalign tool aligns any short protein or nucleotide sequence with database sequences. New options provide additional ways to calculate and display frequency by position, and allow the user to include the surrounding region in the alignment. *08 February 2012*

Questions or comments? Contact us at seq-info@lanl.gov

---

## HIV sequence database

| DATABASES | SEARCH | ALIGNMENTS | TOOLS | PUBLICATIONS | GUIDES | | Search Site |

Search DB
Advanced Search
Intra-patient Search
Next-gen Sequences
Geography

### HIV Sequence Database

**Programs and Tools**

**Search Interface** retrieves HIV and SIV sequences, which can then be aligned and used to build trees

**Geography Search Interface** retrieves HIV sequences based on geographical distribution

**Tools for working with sequences** lists all our online tools, organized by function

**Alignments**

**HIV Premade Alignments** Includes Consensus and Ancestral Sequences, Subtype Reference Alignments, and Complete Alignments

**Information**

**HIV Sequence Compendium** print or order our annual publication

**Tutorials and other information** unpublished web-based content

**Links** to other HIV/AIDS tools and information

**About this website**

**FAQ** general information about this website

**Site Statistics** usage information for www.hiv.lanl.gov

**How to Cite this Database**

**Editorial Board**

### News:    Archived News ▶

**New Features for HIV BLAST**
HIV BLAST has new features. It now allows the user to find best matches among only subtyped sequences, or sequences of a specific subtype. It allows the resulting sequences to be downloaded fully aligned. *01 March 2012*

**New Option for N-GlycoSite**
The N-GlycoSite tool predicts N-linked glycosylation sites in amino acid sequences. A new option allows the user to exclude sites with a second-position proline, which is disfavored for N-linked glycosylation. *29 February 2012*

**HIV Antibody Search Results More Specific**
The antibody search interface in the HIV Immunology database is now more specific. Searches from the Author, Keyword, and Note fields now display only those notes and references that relate directly to the search. The user may still opt to display all, if desired *09 February 2012*

**New Options for Quickalign**
The Quickalign tool aligns any short protein or nucleotide sequence with database sequences. New options provide additional ways to calculate and display frequency by position, and allow the user to include the surrounding region in the alignment. *08 February 2012*

last modified: Tue Jan 26 10:10 2010

Questions or comments? Contact us at seq-info@lanl.gov.

DEPT OF HEALTH & HUMAN SERVICES    NATIONAL INSTITUTES OF HEALTH

# Search Interface

- **Help**
  - □ Tips at the top of the page are often overlooked
    - ■ Ranges, operators, wildcards, logical groupings
  - □ Mouse-over provides brief descriptions; click field names for details in Help file
- **Searches**
  - □ Searches are case-insensitive
  - □ Records are searchable through sequence, patient, genomic region, or publication information and can be matched to the genomic region of a user-provided alignment
  - □ First seven fields will appear in search results page by default
  - □ A "*" in a textbox will cause that field to be included in the results page
  - □ Patient information (Infection year, Infection country) is different than sequence information (Sampling year and Sampling country)
  - □ Problematic sequence filters (hypermutation, frequent ambiguities, potential contamination)
- **Analysis**
  - □ Build a tree with user alignment, search results and subtype reference sequences combined
- **Results**
  - □ Can download aligned or unaligned sequences
  - □ Alignments are based on multiple pair wise alignments – alignments are good, but need hand editing for an optimal alignment
  - □ Select all or a subset of sequences for download
  - □ Sequences can be re-ordered by clicking on fields at the top of the page

Los Alamos
NATIONAL LABORATORY

---

HIV sequence database

| DATABASES | SEARCH | ALIGNMENTS | TOOLS | PUBLICATIONS | GUIDES | | Search Site |

**Sequence Search Interface**

**Tips**
- Click or mouse over the field name for specific tips
- The *italicized fields* are listed in output by default
- To list fields that are not listed by default or included in the search, put an asterisk (*) in the input box
- Use the + and - to see more or fewer search fields
- For other details about each field, see Help or Data Dictionary

Last GenBank update: 2012-02-08
Advanced Search

⊟ **Sequence Information**

| *Accession number* | | | Virus | HIV-1 |
| *Sequence name* | | | Subtype | Any subtype / No subtype / A / A1 / A2 / B |
| *Sequence length* | | | | |
| exact ☑ *Sampling year* | | | | |
| *Sampling country* | BR | | | ☐ Include recombinants |

We will search for country = Brazil (BR)

⊞ More sequence information

⊟ **Find all sequences for a specific gene or region (HIV-1 and SIVcpz)**

| Genomic region | Any / complete genome / 5' LTR / 5' LTR R / 5' LTR U3 / 5' LTR U5 / TAR | Or define start | and end |
| | | ☐ Include fragments of minimum length | 100 |

We will search for complete genomes.

⊞ Combine database sequences with your own sequence alignment (HIV-1 and SIVcpz)
⊞ Publication Information
⊞ Patient Information
⊞ Geographical Information
⊟ **Output**

| | ☐ Include problematic sequences | % of non-ACGT | |
| List | 100 | records per page | Show results selected ☐ | Show SQL ☐ |
| Advanced Search | | | Search | Reset |

last modified: Wed Dec 7 14:05 2011

os Alamos
ONAL LABORATORY

Questions or comments? Contact us at seq-info@lanl.gov.

# Results for HIV-1 complete genomes from Brazil



Choose "One sequence/patient" to remove very similar sequences (only available if a region is selected)



Select a few sequences and make tree, allows us to add a reference set to our data and align them

# TreeMaker tool



**HIV sequence database**

DATABASES  SEARCH  ALIGNMENTS  TOOLS  PUBLICATIONS  GUIDES  | Search Sit |

Choice of outgroup influences the the tree. In general, choose next closest sequences to the "ingroup". In this case our Brazilian sequences are all HIV-1 M group.

These settings minimally influence relative branch lengths, but rarely alter the tree topology.

**Model parameters**
Distance model  Felsenstein 1984 (F84)
Gap handling  ◉ strip gaps before analysis  ○ treat as missing
Site rates  ◉ Equal  ○ Gamma Shape

**Reference sequences (TATCDS)**
○ All
◉ A-K
○ N, O, CPZ, CRFs
○ Menu select only

A1.KE.1994.Q23_17.AF004885
A1.SE.1994.SE7253.AF069670
A1.UG.1985.U455_U455A.M62320
A1.UG.1992.92UG037.U51190
A1.UG.1998.98UG57136.AF484509

**Outgroup**
O.BE.1987.ANT70.L20587
O.CM.1991.MVP5180.L20571
O.CM.1998.98CMU2901.AY169812
O.SN.1999.99SE-MP1299.AJ302646
O.SN.1999.99SE-MP1300.AJ302647

Reference sequences

B.BR.1990.BZ167.AB485641
B.BR.1990.BZ167.AB485642
F1.BR.1990.BZ163.AB485656
F1.BR.1990.BZ163.AB485657
F1.BR.1993.93BR020_1.AF005494

Our Brazilian sequences

Database sequences

Optional mailback, and tree title

**Results link**
Email a link to the results to this address  with job title  Brazil complete genomes

Submit  Reset

---



ATV java-based view for quick look, cannot save/print

**HIV sequence**

DATABASES  SEARCH  ALIGNMENTS  TOOLS  PUBLICATIONS  GUIDES

**Download Your Tree Results**

This tree contains 59 sequences and is 7897 characters long, including insertions.

**Phenogram:**
- View Tree in ATV (a Java-based phylogenetic tree viewer)
- Download Phenogram (pdf)
- View Phenogram (png)

**Radial:**
- Download radial (unrooted) tree (pdf)
- View radial (unrooted) tree (png)

**Alignment used for tree building**
- Download fasta alignment (before gapstripping)
- Download fasta alignment in tree order (before gapstripping)
- Download fasta alignment (after gapstripping)
- Download Newick Tree File

last modified: Thu May 7 07:39 2009

Questions or comments? Contact us at seq-info@lanl.gov.

Obtaining your sequences of interest and having them aligned to a good reference set was the whole point of this. The tree was just a first check on data and alignment quality.

Save alignment, use BioEdit or SeAl to view/adjust.

Save alignment, use BioEdit or SeAl to view/adjust.

Send alignment to GeneCutter or HIV-Align first, is usually best.

http://www.hiv.lanl.gov/content/sequence/GENE_CUTTER/cutter.html

Brazil Genomes Plus Subtype Reference Set, as downloaded



Quick Alignment from Search Interface has many "broken codons"

Send the file through GeneCutter alignment tool to "Codon Align"

New search: all complete genomes; then look at geographic and subtype distribution of the sequences



# Geography output

Each continent's pie chart is clickable to "zoom in" on that continent.

Likewise for each country once you are zoomed in to the continent level.

Most complete genomes in the HIV database are subtype B. But subtype C is more prevalent in human infections. Beware of this type of sampling bias.

New search: all sequences from Brazil. Then look at the distribution of the sequences over the genome

# Tools

- **Analysis and Quality Control**
  - □ **HIV BLAST** finds sequences similar to yours in the HIV database.
  - □ **N-Glycosite** finds potential N-linked glycosylation sites.
  - □ **RIP 3.0** (Recombinant Identification Program) detects HIV-1 subtypes and recombination.
- **Alignment and sequence manipulation**
  - □ **HIValign** uses our HMM alignment models to align your sequences.
  - □ **Gapstreeze** removes columns with more than a given % of gaps.
  - □ **EpimDupes** Given an alignment or set of unaligned nucleotide or protein sequences, this tool compares the sequences and eliminates any duplicates.
- **Phylogenetics**
  - □ **TreeMaker** generates a neighbor-joining phylogenetic tree.
  - □ **PhyML** generates a maximum likelihood phylogenetic tree.
  - □ **TreeRate** finds the phylogenetic root of a tree and calculates evolutionary rate.
- **Format and display**
  - □ **Protein Feature Accent** provides an interactive 3-D graphic of HIV proteins; the user can map a sequence feature (a short functional domain, epitope, or amino acid) and see where it occurs spatially in the 3D structure.
  - □ **Highlighter** highlights mismatches, matches, transition and transversion mutations, and silent and non-silent mutations in an alignment of nucleotide sequences.
  - □ **SeqPublish** makes alignment publication-ready.
  - □ **Recombinant HIV drawing tool** highlights regions of the genome on a graphically representation

Los Alamos
NATIONAL LABORATORY

---

## The HIV database sequence analysis tool set



Click top level to link to full page of tools

Los Alamos
NATIONAL LABORATORY

## HIV Database Tools

(alphabetical order within category)

*For detailed descriptions, mouse over the links.*

### Analysis and Quality Control

Entropy quantifies positional variation in an alignment using Shannon Entropy

HIV BLAST finds sequences similar to yours in the HIV database

Hypermut detects hypermutation

jpHMM at GOBICS detects subtype recombination in HIV-1; hosted at GOBICS as a collaboration between the Department of Bioinformatics, University of Göttingen and the Los Alamos HIV Sequence Database

N-Glycosite finds potential N-linked glycosylation sites

PCOORD multidimensional analysis of sequence variation

Quality Control runs several tools to allow quick QC analysis of HIV-1 sequences; optional step prepares sequence submission for GenBank

RIP (Recombinant Identification Program) detects HIV-1 subtypes and recombination

SNAP calculates synonymous/non-synonymous substitution rates

SUDI Subtyping plots the distance of your sequence to established subtypes

VESPA (Viral Epidemiology Signature Pattern Analysis) detects residues with different frequencies in two sequence sets

### Alignment and sequence manipulation

Codon Alignment takes a nucleotide alignment and returns a codon alignment and translation

Consensus Maker computes a customizable consensus

ElimDupes compares the sequences within an alignment and eliminates any duplicates

Gap Strip/Squeeze removes columns with more than a given % of gaps

Gene Cutter clips genes from a nucleotide alignment, codon-aligns, and translates

HIValign uses our HMM alignment models to align your sequences

### Phylogenetics

Branchlength calculates branch lengths between internal and end nodes

FindModel finds which evolutionary model best fits your sequences

PhyloPlace reports phylogenetic relatedness of an HIV-1 sequence with reference sequences

PhyML generates much better trees than our simple TreeMaker tool

Poisson-Fitter estimates time since MRCA and star-phylogeny. For use with acute (low diversity) samples.

TreeMaker generates a quick-and-dirty phylogenetic tree

TreeRate finds the phylogenetic root of a tree and calculates evolutionary rate

### Immunology

ELF (Epitope Location Finder) identifies known and potential epitopes within peptides

Epilign (QuickAlign) aligns a protein sequence (e.g., epitope) to the appropriate protein alignment

Heatmap displays a table of numbers by using colors to represent the numerical values

Hepitope identifies potential epitopes based on HLA frequencies

Mosaic Vaccine Tool Suite designs and assesses polyvalent protein sequences for T-cell vaccines

Motif Scan finds HLA anchor motifs in protein sequences for specified HLA serotypes, genotypes or supertypes

PeptGen generates overlapping peptides from a protein sequence

### Database search interfaces

ADRA Antiviral Drug Resistance Analysis, a resistance mutation database

Advanced Search creates a custom search interface

Tools are organized in groups by function/purpose.

Most tools have explanation pages, and sample data sets.

Many tools were inspired by user comments, please ask for more.

Los Alamos
NATIONAL LABORATORY

---

SynchAlign aligns overlapping alignments to one another

QuickAlign (formerly Epilign and Primalign) aligns a nucleotide or protein sequence (e.g., primer or epitope) to the appropriate genome alignment

Codon Alignment takes a nucleotide alignment and returns a codon alignment and translation

ElimDupes compares the sequences within an alignment and eliminates any duplicates

Pixel generates a PNG image of an alignment using 1 or more colored pixel(s) for each residue

PepMap can be used to map epitopes, functional domains, or any protein region of interest

### Format and display

Protein Feature Accent provides an interactive 3-D graphic of HIV proteins; can map a sequence feature (a short functional domain, epitope, or amino acid) and see it spatially

Format Converter converts between alignment formats

SeqPublish makes publication-ready alignments

Highlighter highlights mismatches, matches, transitions and transversion mutations and silent and non-silent mutations in an alignment of nucleotide sequences

Recombinant HIV-1 Drawing Tool creates a graphical representation of your HIV-1 intersubtype recombinant

Protein Structure Analysis provides a visualization tool for protein sequence properties

Advanced Search creates a custom search interface

Geography shows the geographic distribution of sequences in the database

CTL/CD8+ Search searches for CD8+ epitopes by protein, immunogen, HLA, author, keywords

T-Helper/CD4+ Search search for CD4+ epitopes by protein, immunogen, HLA, author, keywords

Antibodies search for HIV antibodies by protein, immunogen, AB type, isotype, author, keywords

Vaccine Trials Database finds past vaccine trials and their results

ADRA Antiviral Drug Resistance Analysis, a resistance mutation database

### Other tools

HDent and HDdist perform analysis of heteroduplex mobility shifts

ODprep and ODfit calculate antibody titers based on concentration and optical density data

### External tools

External tools lists tools and programs on other websites

We tend to list only tools of great use in HIV research. Many of these tools are essential, such as either BioEdit or SeAl for alignment viewing and correction.

http://www.hiv.lanl.gov/content/sequence/HIV/HIVTools.html

Los Alamos
NATIONAL LABORATORY

# Pre-Built Sequence alignments

- Originally based on iterations of manual and HMM alignments
- Yearly updates using HMM and manual corrections
- Alignments are in reading frame (codon aligned)
- Contain non-redundant data (one sequence per patient)
- Compendium alignments show fewer sequences than web version
- Reference alignments contain up to four representatives of each subtype.  One of each CRF.
- Protein alignments may contain frameshift compensations
- Subtype consensus with ties resolved, as well as maximum likelihood ancestors, are available for reagent production
- Special interest alignments are being added
  - Sequence sets of particular research interest
  - Suggestions welcome to tkl@lanl.gov



---



Alll(complete) = one per patient, all sequences for which we have a complete genome, or a complete gene.

Subtype Reference = 4 representatives of each subtype, plus one of each Circulating intersubtype recombinant form (CRF) of the M group, plus 4 O group, N group,  P group and SIV-CPZ

Consensus/Ancestral computed from master alignment periodically.

HIV-2/SIV-SMM and primate lentivirus alignments also available here.

# SIV/PLV Alignments

- Any non-human lentivirus is a SIV (or primate lentivirus), not just the SIV-SMM/SIV-MAC group from Sooty mangabeys.
- HIV-1s (M, N, O and P groups) are related to the SIV-CPZs from the chimps (*P. t. troglodytes)* and SIV-GORs from gorillas. We describe these alignments as HIV-1/CPZ.
- HIV-2s and SIV-MACs are related to SIV-SMMs from Sooty mangabeys. We describe these alignments as HIV-2/SMM.
- Dozens of other diverse non-human primates, such as African green monkeys, carry species-specific SIVs.
- Alignments of the diverse SIVs, plus HIVs, can help to identify highly conserved codons and other features. We describe these alignments as "other SIV" or HIV-1/HIV-2/SIV.

**Los Alamos**
NATIONAL LABORATORY

---

# Primate Lentiviruses

**Alignments: http://www.hiv.lanl.gov/content/hiv-db/ALIGN_CURRENT/ALIGN-INDEX.html**



*P. t. troglodytes*

X

**Positive Chimps
HIV-1 M, N, O**

*P. t. schweinfurthii*

**Van Heuverswyn, Nature 2006
Keele, *Science* 2006
Corbet, *J. Virol* 2000
Foley, HIV database**

**Los Alamos**
NATIONAL LABORATORY

# Gene Cutter

- Unconventional Alignment/Homology program
- "Cuts out" specified genes and proteins from sets of DNA sequences
  - ☐ Aligns to HXB2 via HMMer (or to SIV-Mac239 for HIV-2 and SIV-SMM)
  - ☐ Splits input sequences into genes, if desired
  - ☐ Aligns DNA sequences by codon, and translates them (including interpretation of IUPAC codes such as R for purine)
- Useful for processing new sequence data
  - ☐ annotating full length genomes
  - ☐ pulling out regions of interest from raw sequence data
- For each gene/region, maintains a list of anomalies
  - ☐ stop codons
  - ☐ codons containing multi-state characters
  - ☐ codons containing indels
- Input sequences may be aligned or unaligned
- Results may be better if the HXB2 sequence is included as a reference in your input file

Los Alamos
NATIONAL LABORATORY

---

# GeneCutter

### Gene Cutter: Sequence Alignment and Protein Extraction

**Purpose:** Gene Cutter is a sequence alignment and protein extraction tool. It can be used for any set of nucleotide sequences for HIV-1, HIV-2 or SIV.

Gene Cutter can:

- align your nucleotide sequences (if they aren't already aligned)
- clip pre-defined coding regions from a nucleotide alignment
- codon-align the coding regions
- generate nucleotide and protein alignments of the cut regions

**Details:** The reference sequence used by this tool is HXB2(Accession #K03455) for HIV-1 or SMM239(Accession #M33262) for HIV-2 or SIV. Gene coordinates are based on these reference sequences. This version of Gene Cutter doesn't require a reference sequence to be included in your input nucleotide alignment. Gene Cutter will also accept **unaligned** sequence sets. Gene Cutter uses Hmmer with a training set of the full-length genome alignment and will give a better multiple alignment than many computationally-based alignment programs. Misalignments at the ends of a coding region may result in a few amino acids/bases not appearing in the output for that coding region.

In some sequences, an insertion will be compensated within a short distance by a deletion, or vice versa. As these frameshifts may not inactivate the protein, if a compensating mutation is within 5 amino acids of an initial frameshift, the shifted reading frame is left intact. Otherwise, the frame shift is marked with the hash symbol (#), and the translation is continued in the correct reading frame beyond the offending codon. Stop codons are marked by a dollar sign ($).

**The best results** will be obtained if you submit an alignment that has been hand-aligned and contains the correct reference sequence. For more information, see Gene Cutter Explanation.

**Input**

Select the organism  [ HIV1 (HXB2) ▼ ]
Paste your sequences
[Sample Input]

Or upload your file:  /Users/btf/Desktop/Outputs/OurData-PlusRefset.FASTA   [Browse...]
Check box if appropriate  ☐ Sequences are unaligned

**Options**

Region(s) to align and extract  [ Env CDS ▼ ]
☐ Insert HXB2(Accession #K03455) for HIV-1 or SMM239(Accession #M33262) for HIV-2 or SIV into the result set
☐ Remove HXB2(Accession #K03455) for HIV-1 or SMM239(Accession #M33262) for HIV-2 or SIV from the result set
☑ Codon align the region

**Translation options**

○ Codons containing an IUPAC character are shown as "X".
○ Codons containing an IUPAC character in a silent position are translated; others are shown as "X".
○ Codons containing an IUPAC character are translated.
◉ Do not translate to amino acids
Note: codons containing "-" are always translated to either "-" (gap) or "#" (partial codon)

[ Submit ]  [ Reset ]

**Please be patient.** Your input file must download to our server, where the actual work is performed. This can take several

Input is our data plus the "reference Set" and any other sequences we chose to add from the search interface.
Input: GeneCutterInput.FASTA
Output: GeneCutterOutputAll.FASTA

For this exercise, we want the Env gene, codon aligned, but not translated to proteins.

Output: GeneCutterOutputEnv.FASTA

Los Alamos
NATIONAL LABORATORY

# GeneCutter Results

## Gene Cutter Mailback Form

Please enter the email address to send the results set: [                    ]

[ Submit email address ]

- Results are stored on our server
  - An HTML link is e-mailed to the user when the run is complete
  - For this workshop, we will provide example.

Los Alamos
NATIONAL LABORATORY

---

# GeneCutter Result

Result saved in Outputs folder
Alignments viewed with Pixel
http://www.hiv.lanl.gov/content/sequence/pixel/pixel.html

Our data aligned to reference set by search tool:
GeneCutterInput.FASTA
(output of search and tree build was input to GeneCutter)



Our data aligned to reference set by GeneCutter:
Outputs: GeneCutterOutputENV.FASTA



Can also be viewed with BioEdit, Se-Al or other multiple sequence alignment editors.

Los Alamos
NATIONAL LABORATORY

# Treemaker

Check for phylogenetic relatives:

- TreeMaker produces a Neighbor Joining tree for a quick comparison

- TreeMaker uses PAUP* for its calculations; a few model options are available

- Reference sequences can be included, and are aligned to the input automatically

- Trees are displayed using PHYLIP and ATV

- The alignment used for the tree can also be downloaded

- A Phyml interface is also available
  http://www.hiv.lanl.gov/content/sequence/PHYML/interface.html

---

http://www.hiv.lanl.gov/components/sequence/HIV/treemaker/treemaker.html

| DATABASES | SEARCH | ALIGNMENTS | TOOLS | PUBLICATIONS | GUIDES | | Search Site |

### Neighbor TreeMaker

**Purpose:** This tool takes a nucleotide sequence alignment, converts it to NEXUS format, and uses PAUP to generate a tree, which is displayed using the PHYLIP programs Drawgram or Drawtree.

**Details:** After sequence input, the next page will give additional options. Gaps can be treated as missing or stripped. The user can choose from various distance models and select the outgroup sequence. A version of the input alignment in which the sequences have been reordered to match the order in the tree may be downloaded. Trees are calculated using the neighbor-joining method. You can use FindModel to decide what evolutionary model best fits your data.

**Disclaimer:** This interface only offers very basic, 'quick-and-dirty' phylogenetic analysis. More in-depth analysis is usually needed. For more information see the Tree Tutorial.

**Input**

Paste alignment here
[Sample Input]

Paste or type a DNA **alignment** here.

or upload your file [          ] Browse...

**OR** upload an alignment file here.

**Tree parameters**

Include reference sequences (HIV-1/CPZ only) ☐

Submit   Reset

http://tree.bio.ed.ac.uk/software/figtree/

# HIV/SIV Sequence Locator Tool

- Instantly computes position numbers of DNA or protein fragments relative to a reference strain (HXB2r for HIV-1, SMM239 for SIV)
  - □ Such numbers, often included in the literature, are frequently incorrect
- Shows the location of the sequence on an HIV map
- Presents protein translations of DNA sequences
- Can be used for input into the search interface, to align a new sequence you have generated with the database set
- Can also retrieve reference sequences
  - □ by coordinates (range of base or amino-acid positions)
  - □ by single position (retrieves flanking sequences)

http://www.hiv.lanl.gov/content/sequence/LOCATE/locate.html

**HIV Sequence Locator Tool**

**Purpose:** This tool has several purposes. It can find the start and end coordinates (relative to the reference strain HXB2) of your input sequence(s) and show which genes or proteins it covers, along with a graphical view of the location of your sequence(s) relative to the reference sequence. The tool will display both the nucleotide sequence and protein translation of your input as it aligns to HXB2. It will also check the reverse complement of your input sequence, and report the orientation with the best match. Another use is to retrieve a section of the HXB2 reference sequence based on its coordinates.

**How to use:** To find the coordinates for your sequence, either upload or paste your sequence (any format) in the box below, or (for database sequences only) enter GenBank accession numbers. To retrieve the HXB2 sequence for a set of coordinates (see HIV coordinate map), enter the coordinates and choose the region. To retrieve the entire gene or protein, enter coordinate values of "1" and "end". To retrieve a single nucleotide or range with its surrounding 42-nucleotide sequence, enter the single coordinate in the "from" field and check the box. For more details, see Sequence Locator Explanation.

**Useful Links:**
HXB2 numbering | SIVmm239 numbering (review articles)
HXB2 spreadsheet | SIVmm239 spreadsheet (spreadsheets with base-by-base annotation)

Paste or type a DNA or protein sequence here.

OR enter numeric coordiantes here.

---

**Sequence Locator: "find my sequence"**



Result for Sample Input DNA Query sequence

Location in genome mapped in red.

Numeric coordinates useful for entry on search form

DNA and protein sequence displayed

## Sequence Locator: "Retrieve from coordinates"

Table of genomic regions touched by query sequence. Query protein translation in blue.

| CDS | NA position relative to CDS start in HXB2 | NA position relative to query sequence start | NA position relative to HXB2 genome start | AA position relative to protein start in HXB2 |
|---|---|---|---|---|
| Gag | 352 -> 483 | 1 -> 132 | 1141 -> 1272 | 118 -> 161 |
| AAADTGHSNQVSQNYPIVQNIQGQMVHQAISPRTLNAWVKVVEE | | | | |
| p17 | 352 -> 396 | 1 -> 45 | 1141 -> 1185 | 118 -> 132 |
| AAADTGHSNQVSQNY | | | | |
| p24 | 1 -> 87 | 46 -> 132 | 1186 -> 1272 | 1 -> 29 |
| PIVQNIQGQMVHQAISPRTLNAWVKVVEE | | | | |

Sequence below includes up to 42 bases of context surrounding query sequence.

| Reference Strain | Type | Region | Start | End |
|---|---|---|---|---|
| HXB2 | nuc | complete | 1141 | 1272 |

Retrieved Sequence:

GCAGCAGCTGACACAGGACACAGCAATCAGGTCAGCCAAAATTACCCTATAGTGCAGAACATCCAGGGGCAAATGGTACATCAGGCCATATCACCTAGAACTTTAAATGCATGGGTAAAAGTAGTAGAAGAG

Organism: HIV

Los Alamos
NATIONAL LABORATORY

---

# HIV Genome Browser:

- Dreamed of by Christian Brander and designed by Shihai Feng, with the help from Jennifer Macke, Brian Foley, Jim Szinger, Karina Yusim

- A customization of Jbrowse Genome Browser, built to incorporate many sources of information from the LANL HIV Sequence and Immunology databases.

- A one-stop source of information about HIV genome and immunological data. It retrieves the vast and diverse information available at HIV Immunology database and allows to look at the whole HIV genome as well as zoom in to a region of interest and see all information we have in the database about this region
    - HXB2 gene map, HXB2 sub-protein map, Mac239 map
    - Overlapping epitopes, antibody binding sites
    - HXB2 coding sites of interest (e.g. functional domains, drug resistance sites, motifs, glycosylation sites, etc.)
    - HXB2 LTR sites of interest (RNA structural elements, primer binding sites, etc.)
    - Neutralizing Ab contact residues, signatures and other NAb-associated features
    - HIV sequence variability (Entropy: M group, B clade, C clade)
    - Links to the database annotation, alignments, tools, Pubmed etc.

- DNA- and Protein-level views are available

Los Alamos
NATIONAL LABORATORY

# HIV Genome Browser

**Purpose:** To display graphic views of the HIV genome and proteome, allowing the juxtaposition and exploration of multiple types of data. Details in Help.

## Starting Points

These are just starting examples; within the genome browser, you can move between any of these views.

**Nucleotide-level example views:**

- HIV-1 gene map
- SIV Mac239 gene map
- HIV-1 5' LTR

**Protein-level example views:**

- HIV-1 Env: CTL epitopes + entropy
- HIV-1 Pol: drug resistance sites + entropy

## Quick tips

- Use mouseovers! There are many mouseovers to guide you.
- Use click and right-click! Every feature has additional information and analysis available via click or right-click. If your mouse doesn't have right-click, use Ctrl-click.
- Zoom! There are several ways to zoom in and out. Some features can only be seen when zoomed-in or zoomed-out.
- For details about this interface, see HIV Genome Browser Help.
- Watch the screencast video on the JBrowse website.

---

# HIV Genome Browser: Nucleotide view

## HIV genome browser: more possibilities ?

- Data of how heavily sequenced each genome region is (we are getting questions sometimes why some regions don't return a lot of sequences on the sequence search interface)

- Show subtype consensus sequences

- CTL Epitope variants (we currently have a database of ~3000 CTL epitope variant records and started Helper epitope variants)

- Categorize heavily loaded tracks. For example, provide separate tracks for Drug resistance, CD4 contact residues, Ab contact residues, Glycosylation sites etc

- Links to structure

- Suggestions ?

Los Alamos
NATIONAL LABORATORY

---

# Hypermutation

**Hypermut 2.0**

**Analysis & Detection of APOBEC-induced Hypermutation**

**Purpose:** This interface takes a nucleotide alignment and documents the nature and context of nucleotide substitutions in a sequence population relative to a reference sequence.

**Details:** The first sequence in the input alignment will be used as the reference sequence, and each of the other sequences will be used as a query sequence. Please choose the reference sequence carefully. For example, for an intrapatient set, the reference should probably be the most common form in the first sampled time point; for a set of unrelated sequences, the reference should probably be the consensus sequence for the appropriate subtype. Before using, please read:

- Hypermut Explanation
- Hypermut 2.0 Details

**References:** Please reference these articles when using Hypermut:

- Rose, PP and Korber, BT. 2000. Detecting hypermutations in viral sequences with an emphasis on G -> A hypermutation. *Bioinformatics* **16**(4): 400-401.
- Bruno, WJ, Abfalterer, WP, Foley, BT, Leitner, TK and Korber, BT. Detection of hypermutation in HIV sequences using two context positions and avoiding nucleotide content effects. Manuscript submitted.

**Input**

Indicate sequence format of input | FASTA
Note: Sequences must be aligned, in-frame if possible, and of equal length.
Paste alignment here
>Seq1
CAACTGCTGTTAAATGGCAGTCTAGCAGAAGAAGAGGTAGTAATTA
GATCTGAAAATTTCACCAATAATG
CTAAAATCATAATAGTACAGTTGAATGAATCTGTAAAAATTGATTG
TATAAGACCCAACAACAATACAAG
AAAAAGTATACATATCGGACCAGGGAGAGCATTTTACACAACAGGA

Or upload alignment file | Choose File | no file selected
Restrict analysis to subregion of alignment from | bp to | bp (optional)

**Hypermut 2.0 Customized Options**

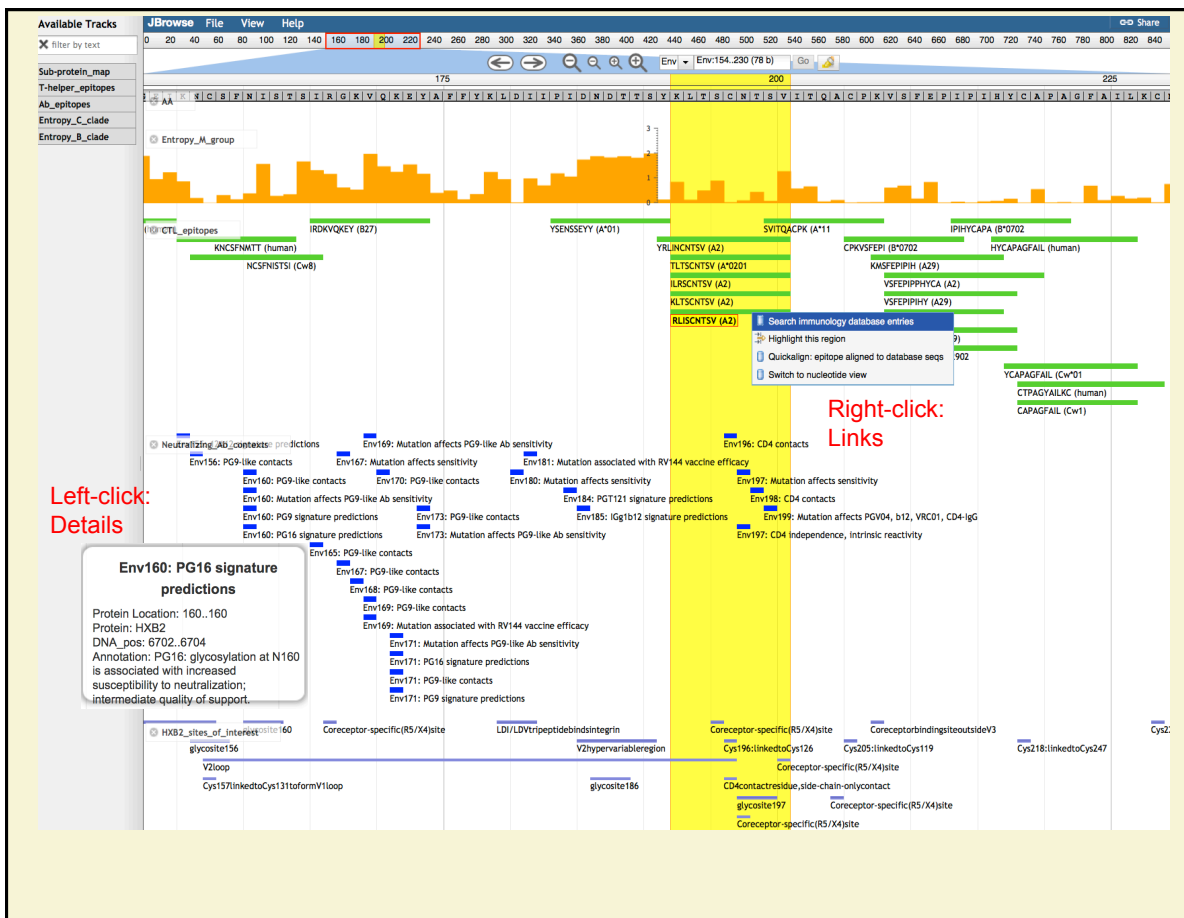These options apply only to Hypermut 2.0 analysis, and have no effect on the Original Hypermut output. For typical analyses of APOBEC-induced hypermutation in HIV, these options should be left in their default settings.

Mutation
Upstream context: ↓ Downstream context: Enforce context:

Customize Hypermut pattern: | G | RD | ○ On reference sequence
↓ | ○ On both sequences
A | ◉ On query sequence

Customize control pattern: | G | YN|RC
↓
A

**Output**

Analyses to perform: ◉ Both ○ Original Hypermut ○ Hypermut 2.0

Run | Reset

- Detects APOBEC related A->G hypermutation as default

- Can be adapted to detect any fuzzy motif in relation to a control pattern

Los Alamos
NATIONAL LABORATORY

**Hypermut results**

**Hypermut 2.0**

Your pattern definitions are as follows. Where there is no pattern (i.e., just '...') all sequences will match.

| Pattern | Upstream | From →To | Downstream |
|---------|----------|----------|------------|
| 'Mut' | ... | G →A | RD ... |
| 'Control' | ... | G →A | YN|RC ... |

**Results**

'Potential Mut' or 'Potential Control' means a match to the corresponding Upstream, From, and Downstream patterns above, while an actual 'Mut' matches those and the To pattern as well. We consider a P-value less than 0.05 to indicate a hypermutant when using the default patterns.

| Sequence: (Select for graphing) | Muts: (Match Sites) | Out of: (Potential Mut Sites) | Controls: (Control Muts) | Out of: (Potential Controls) | Rate Ratio: (A/B)/(C/D) | Fisher Exact P-value: (=P(Muts,Poten.Muts-Muts, Cntrls,Poten.Cntrls-Cntrls)) |
|---|---|---|---|---|---|---|
| ☑ Seq2 | 0 | 71 | 0 | 54 | undef | 1 |
| ☑ Seq5 | 4 | 69 | 1 | 52 | 3.01 | 0.282669 |
| ☑ Seq7 | 26 | 71 | 1 | 54 | 19.77 | 5.35061e-07 |
| ☑ Seq14 | 48 | 71 | 9 | 54 | 4.06 | 8.26961e-09 |

High ratio of G -> A vs. A -> G indicates hypermutation

**View Sites Along Sequence**

Type of graph:
- ⦿ Locations of Matches
- ○ Cumulative Matches (try me!)

(Graph Matches) (opens in a new window)

**Optional Controls:**
Show region: From [ ] to [ ]
Graph Title: Hypermut Custom Analysis
Access xmgrace compatible datafile.

Cumulative mutation Graph is useful

**Original Hypermut Output**

The input file has 5 sequence(s)
Sequence Length: 645
Compared to SEQ1, 264 As, 126 Gs, 92 Cs, 163 Ts

| | Sequence names | Ratio | #diffs | perc_Gs | #A->G | #G->A | Dinuc Context: GG GA GC GT | OBSERVED CHANGES |
|---|---|---|---|---|---|---|---|---|
| ○ ○ | SEQ2 | 0/0 | 1 | 0.00 | 0 | 0 | 0 0 0 0 | TC |
| ○ ○ | SEQ5 | 5/2 | 33 | 3.97 | 2 | 5 | 2 3 0 0 | GA TA GA CT CA GA AT CT AG A- T- A- G- T- |

*Hypermut Custom Analysis — Mutation: G --> A, Context: _RD, ControlContext: _YNIRC*

---

# Highlighter

- Highlights mutations relative to a reference strain, particularly useful for intra-patient analyses.
- Highlights:
  - ☐ syn/non-syn
  - ☐ transition/transversion
  - ☐ Apobec motifs
- Sorts on similarity
- Visualize recombination of closely related sequences

Sequences compared to master

Nonrandom distribution of mutations evident.

Sample Set is from a possible dual Infection, with intra-subtype recombinants evident.

---

# Protein Feature Accent

- Highlights region of interest in an HIV structure
- You can upload a PDB structure, or use one of our annotated Env structures
- You can upload your own alignment and get an entropy map

http://www.hiv.lanl.gov/content/sequence/PROTVIS/html/protvis.html

**HIV sequence database**

DATABASES    SEARCH    ALIGNMENTS    TOOLS    PUBLICATIONS    GUIDES    [          ]    Search Site

## Protein Feature Accent

**This is a beta version!**
Some capabilities are not fully implemented, and there may be bugs or other problems. Please use with care and a sense of humor.
This tool requires that Java be installed on your computer.

**Purpose:** The Protein Feature Accent tool is a quick way to map protein sequence features (for example a short functional domain or an epitope) from a sequence directly on to an interactive graphic of the corresponding 3-D structure of the protein.

**How to use:** The tool needs only to be directed to use a particular protein structure file in PDB format. Uploading a sequence is not required; the sequence associated with the chosen structure will always be displayed. Any sequences you do provide will be analyzed (for entropy, etc.), aligned with the structure sequence, and displayed.
If you prefer, you may upload a PDB file for a structure you wish to use instead of those available here.
Click here for a list of all the structures available.

New features:
- Predicted N-linked glycosylation site highlighting
- User-supplied alignment entropy color scheme
- PDB file upload option

We are in the process of adding additional features to the tool.

Select a protein structure:  **gp120**
[switch to full structure list]
       1G9M: HIV-1 HXBC2 GP120 ENVELOPE GLYCOPROTEIN COMPLEXED WITH CD4 A . . .
      2B4C: HIV-1 JR-FL GP120 CORE PROTEIN CONTAINING THE THIRD VARIABLE . . .
HIV ⦿ SIV ◯  2NY7: HIV-1 GP120 ENVELOPE GLYCOPROTEIN COMPLEXED WITH THE BROADLY . . .
      1RZK: HIV-1 YU2 GP120 ENVELOPE GLYCOPROTEIN COMPLEXED WITH CD4 AND . . .
**gp41**

OR

upload a PDB file:  [          ]  Browse...

You may provide an amino acid sequence alignment (or a single unaligned sequence) below:
Paste your sequence(s) here
```
>B.BR.99.BREPM11931_DQ085869
MRVRETKKNYWQWWRRGMMLLGMLMICSATEQSWVTVYYGVPVWKEASTTLFCASDAKAVETEAHNVWAT
HACVPTDPNPQEVVLENVTENFNMWKNNMVEQMHEDIISLWDQSLKPCVKLTPFCETKMCSNVDNATSDT
NSTNSGWEKMAEEIRNCSFNVTTNIGNKRQKEYALFNKLDVVPIDNTSYTLINCNTSVITQACPKISFEP
IPIHYCTPAGFAILKCNDKKFNGTGPCKNVSTVQCTHGIRPVVSTQLLLNGSLAEEEIVIRSENFTNNAK
TIIVQLNKTVVINCTRPNNNTRKGIHLGPGRTVYATGGIIGNIRQAHCNISGAEWENTLKQIATKLGGQF
KNKTIAFNQSSGGDPEITMHSFNCGGEFFYCNTTQLFNSTWTYTWNRNGNGTNGTITLPCRIKQIINRWQ
```
or upload a sequence file:  [          ]  Browse...

> List of "recommended" PDB entries

> Only a gp120 alignment is provided so far. We hope to add others. You can paste in your own.

---

http://www.hiv.lanl.gov/content/sequence/PROTVIS/html/protvis.html

**Jmol window**  The viewing window below offers Jmol's interactive features, in addition to the control panel at the left.

**Control Panel**

Protein Data Bank structure ID: **2B4C**
JRFL gp120 as complexed with CD4 and Ab X5; has V3 loop, lacks N- & C-terminals and V1/V2 loop

[+]  [↑]  [−]
    [↑]
[⇐] [←] [→] [⇒]
    [↓]
    [⇓]

(re-center)    spin ☐
[select display style ⇕]
[pick color scheme ⇕]
Background: ■■■□□  ■■□□

V1/V2 Loop    V5 Loop
V4 Loop
CD4 bs
V3 Loop

Jmol command script:
[                    ]
(execute)

Jmol_S

(download) this view as a [ jpg ⇕ ] image

> Many display options in JMol are "built in" to this web tool.
> Use the JMol command script box below for other commands.
>
> One of the color schemes is "color by entropy" based on divrsity in the alingment added below.

**Sequence View**
Select residues in the top (PDB file) sequence below to highlight them in the graphic above.
☑ show PDB file annotation   ☑ show reference sequence   ☐ show reference sequence annotation

```
                              320       330       340       350       360       370       380
                               .    .    .    .    .    .    .    .    .    .    .    .    .
      PDB file sequence:-GPGRA-FYTTGEIIGDIRQAHCNISRAKWNDTLKQIVIKLREQFEN-KTIVFNHSSGGDPEIVMHSFNCGGEFFYCNS
                   HXB2:RGPGRA-FVTIGKI-GNMEQAHCNISRAKWNNTLKQIASKLREQFGNNKTIIFKQSSGGDPEIVTHSFNCGGEFFYCNS
                              -GPGRTVYATGGII-GNIRQAHCNLSGAEWENTLKQIATKLGGQF-KNKTIAFNQSSGGDPEITMHSFNCGGEFFYCNS
B.BR.99.BREPM11931_DQ085869:-GPGGTIYATGGII-GNIRQAHCNISGAEWENTLKQIATKLGGQF-KNKTIAFNQSSGGDPEIIMHSFNCGGEFFYCNT
B.BR.99.BREPM11932_DQ085870:-GPGAFYTTGDII-GDIRKAHCNLSKSDWNNALRQVARKLGGQF-KNKTINPTRSSGGDPEIMHSFNCGGEFFYCNS
  B.CA.00.CANA6FULL_AY779552:-GPGAFYTTGEII-GNIRQAHCNLSRAEWNKTLEQIVGKLREQF-GNKTIVFNQSSGGDPEIVTHSFNCGGEFFYCNT
      B.CA.82.82CAN_AY247225:
```

> Selected region gets highlighted in structure

**Highlighting**

# Quality Control Tool

- Built from existing HIV database tools
- GeneCutter
- RIP
- Hypermut
- Neighbor-joining Trees
- Output is an email containing a link to a summary report
  - http://www.hiv.lanl.gov/content/hiv-db/QC/index.html (beta version)

---

# Quality Control Tool

http://www.hiv.lanl.gov/content/sequence/QC/index



Recently added shortcuts to GenBank entry creation tool.

Requires FASTA format sequences, and a comma separated values (CSV) file of annotations, as described on the help page.

http://www.hiv.lanl.gov/content/sequence/QC/field_help.html

Easy to enter in spreadsheet like EXCEL, and then export as CSV format.

# Quality Control Tool

- Summary of results from analysis programs

- Click on each result to obtain full analysis

- Useful for helping to determine subtype, hypermutation, mislabeling of samples

Los Alamos
NATIONAL LABORATORY

# Please leave any comments or suggestions with us:

Los Alamos
NATIONAL LABORATORY